

Elsevier Editorial System(tm) for Combustion
and Flame

Manuscript Draft

Manuscript Number: CNF-D-16-00540R3

Title: Outlier analysis for a silicon nanoparticle population balance
model

Article Type: Accepted Paper

Keywords: Silicon; nanoparticles; population balance; regression
influence diagnostics

Corresponding Author: Professor Markus Kraft,

Corresponding Author's Institution: University of Cambridge

First Author: Sebastian Mosbach, PhD

Order of Authors: Sebastian Mosbach, PhD; William J Menz; Markus Kraft

Manuscript Region of Origin: UNITED KINGDOM

Outlier analysis for a silicon nanoparticle
population balance model

Sebastian Mosbach^a, William J. Menz^a, Markus Kraft^{a,b,*}

^a*Department of Chemical Engineering and Biotechnology, University of Cambridge
New Museums Site, Pembroke Street, Cambridge CB2 3RA, United Kingdom*

^b*School of Chemical and Biomedical Engineering, Nanyang Technological University
62 Nanyang Drive, Singapore 637459, Singapore*

Abstract

We assess the impact of individual experimental observations on a multivariate population balance model for the formation of silicon nanoparticles from the thermal decomposition of silane by means of basic regression influence diagnostics. The nanoparticle model is closely related to one which has been used to simulate soot formation in flames and includes morphological and compositional details which allow representation of primary particles within aggregates, and of coagulation, surface growth, and sintering processes. Predicted particle size distributions are optimised against 19 experiments across ranges of initial temperature, pressure, residence time, and initial silane mass fraction. The influence of each experimental observation on the model parameter estimates is then quantified using the Cook distance and DFBETA measures. Seven model parameters are included in the analysis, with five Arrhenius pre-exponential factors in the gas-phase kinetic rate expressions, and two kinetic rate constants in the population balance model. The analysis

*Corresponding author
Email address: `mk306@cam.ac.uk` (Markus Kraft)

highlights certain experimental conditions and kinetic parameters which warrant closer inspection due to large influence, thus providing clues as to which aspects of the model require improvement. We find the insights provided can be useful for future model development and planning of experiments.

Keywords: Silicon, nanoparticles, population balance, regression influence diagnostics

1. Introduction

Gas-phase synthesis in hot-wall reactors is a common way in which silicon nanoparticles are manufactured. Shock-tubes are another set-up in which especially the early phase of formation of these particles can be studied. Typically, these synthesis processes begin with silane (SiH_4) as a precursor, which is transformed into the eventual nanoparticle product at high temperatures. A variety of models have been proposed to describe this transformation [1]. These models usually contain unknown or low-confidence (kinetic) parameters with large uncertainties associated to them. Systematic parameter estimation techniques can then be employed to arrive at better values for these quantities, based on available experimental data. One of the most elementary parameter estimation methods is least-squares optimisation, *i.e.* minimising the distance between experimental observations and model prediction as measured by a sum-of-squares objective function. The result of such an optimisation is a set of values, called ('best') estimates, for the selected model parameters. Not all experimental data points may equally inform the optimal value of the parameters, though – different parameters may be determined to a varying extent by different observations. In order to assess which ex-

periments are the most relevant in the optimisation, one can conduct what may be called an omission-based regression influence analysis [2]: Firstly, optimise the model against the full data set, and then repeat the optimisation with one of the data points removed, for each of the data points. Based on the difference between the parameter estimates of the full optimisation and the optimisations with an omitted data point, it is then possible to quantify the influence of individual observations on the model overall or on individual parameters. Several such measures have been proposed [3, 4], the most widely-used one being Cook’s distance [5], and applied to detect influential data points, high-leverage points, and statistical outliers [6, 7].

An alternative approach to quantifying influence of experimental observations is uncertainty propagation [8], part of which is concerned with how experimental measurement errors propagate into model parameters and responses. Some of these methods allow calculating the relative contribution of each data point (and its error bar) to the uncertainty in each of the parameters. In particular, the Data Collaboration framework [9] exploits the pairwise consistency of data set units to identify outliers.

Yet another approach, called perturbation of the optimum, has been developed for constrained optimisation [10, p. 34] and unconstrained least-squares optimisation [11], which has found application in chemical kinetics [12, 13, 2]. These methods allow calculating sensitivities of parameter estimates with respect to any other quantity in the objective function (or constraints), including in particular experimental data.

The purpose of this paper is to conduct an omission-based outlier analysis of a selection of experimental data for silicon nanoparticles produced from

1
2
3
4
5
6
7
8
9 a silane precursor in hot-wall flow reactors and shock tubes which are mod-
10 elled using a detailed population balance model. A main aim is to identify
11 those experimental conditions which are the most challenging for the model.
12 We apply a technique established in the field of regression influence diag-
13 nostics to quantify the influence of individual experimental observations on
14 kinetic parameter estimates for this purpose. We determine the influence of
15 the measurements on estimates of some Arrhenius pre-exponential factors in
16 the gas-phase kinetic mechanism as well as the population balance model for
17 the particle phase. Using a threshold for the influence values, specific mea-
18 surements are then highlighted for further analysis, providing further insight
19 into the model and potential improvements, as well as suggestions for future
20 experiments.
21
22
23
24
25
26
27
28
29
30
31

32 33 **2. Background**

34
35 We firstly describe the model, provide some background on omission-
36 based regression influence diagnostics, and how it can be used to identify
37 outliers.
38
39
40
41

42 43 *2.1. Population balance model for silicon nanoparticle formation*

44
45 We briefly summarise the main features of the model here. Full details can
46 be found in [1], and further in [16, 17, 18, 19, 20], noting that a closely related
47 model has been applied to soot formation in flames (see for example [21] and
48 references therein). It consists of two main parts, a gas-phase model, and a
49 particulate phase model.
50
51
52
53
54
55
56
57
58

Table 1: The gas-phase kinetic mechanism. Values in bold correspond to parameters chosen for the influence analysis. Units for the Arrhenius pre-exponential factors are cm, mol, and s.

Idx.	Reaction	A	β [-]	E [kcal/mol]	Ref.
1	$\text{SiH}_4 (+\text{M}) \rightleftharpoons \text{SiH}_2 + \text{H}_2 (+\text{M})$	3.12×10^9	1.7	54.71	[14]
	Low pressure limit:	3.96×10^{12}	0	45.10	[15, 1] ¹
2	$\text{Si}_2\text{H}_6 (+\text{M}) \rightleftharpoons \text{SiH}_4 + \text{SiH}_2 (+\text{M})$	1.81×10^{10}	1.7	50.20	[14]
	Low pressure limit:	5.09×10^{53}	-10.37	56.03	[14]
3	$\text{Si}_2\text{H}_6 (+\text{M}) \rightleftharpoons \text{Si}_2\text{H}_4\text{B} + \text{H}_2 (+\text{M})$	9.09×10^9	1.8	54.20	[14]
	Low pressure limit:	7.79×10^{40}	-7.77	59.02	[14, 1] ²
4	$\text{Si}_3\text{H}_8 (+\text{M}) \rightleftharpoons \text{SiH}_2 + \text{Si}_2\text{H}_6 (+\text{M})$	6.97×10^{12}	1.0	52.68	[14]
	Low pressure limit:	1.73×10^{69}	-15.07	60.49	[14]
5	$\text{Si}_3\text{H}_8 (+\text{M}) \rightleftharpoons \text{Si}_2\text{H}_4\text{B} + \text{SiH}_4 (+\text{M})$	3.73×10^{12}	1.0	50.85	[14]
	Low pressure limit:	4.36×10^{76}	-17.26	59.30	[14]
6	$\text{Si}_2\text{H}_4\text{B} (+\text{M}) \rightleftharpoons \text{Si}_2\text{H}_4\text{A} (+\text{M})$	2.54×10^{13}	-0.2	5.38	[14]
	Low pressure limit:	1.10×10^{33}	-5.76	9.15	[14]
7	$\text{Si}_2\text{H}_4\text{B} + \text{H}_2 \rightleftharpoons \text{SiH}_4 + \text{SiH}_2$	9.41×10^{13}	0	4.09	[14]
	Reverse coefficients:	9.43×10^{10}	1.1	5.79	[14]
8	$\text{Si}_2\text{H}_4\text{B} + \text{SiH}_4 \rightleftharpoons \text{Si}_2\text{H}_6 + \text{SiH}_2$	1.73×10^{14}	0.4	8.90	[14]
	Reverse coefficients:	2.65×10^{15}	0.1	8.47	[14]

¹ A is from [1], β and E are from [15]. ² A is from [1], β and E are from [14].

2.1.1. Gas phase

The gas-phase chemical kinetic reaction mechanism used is a modified version of the one proposed by [14], and is summarised in Table 1. Two isomers of Si_2H_4 are included: silene, *i.e.* H_2SiSiH_2 , denoted by the suffix “A”, and silylene, *i.e.* HSiSiH_3 , denoted by the suffix “B”. The first six reactions are third-body reactions whose pressure-dependence is given in Lindemann fall-off form. More details can be found in [1].

2.1.2. Particulate phase

The particle phase is described by a detailed, high-dimensional population balance model [1] covering aggregate morphology and chemical composition. In this model, each nanoparticle is represented as a list of primary particles, together with a (triangular) matrix, called connectivity matrix, each entry of which represents the common surface area for the corresponding pair of primary particles. For each primary particle, the number of silicon and the number of hydrogen atoms are stored. From this particle representation, beyond elementary properties like mass and chemical composition, several quantities of interest can be derived. These include for example, with some additional assumptions, collision and mobility diameter of aggregates, surface area, and sintering level.

The following processes which create or transform particles, or account for interaction of the particles with the gas phase, are represented in the model:

Inception: Any two molecules of any of the three species SiH_2 , $\text{Si}_2\text{H}_4\text{A}$, and $\text{Si}_2\text{H}_4\text{B}$ can collide to (irreversibly) form a new particle, which is assumed to consist of a single, spherical primary whose diameter follows directly from

1
2
3
4
5
6
7
8
9 its mass, *i.e.* numbers of atoms. The rate at which this happens is as-
10 sumed to be non-zero only if the diameter of the resulting particle exceeds a
11 temperature- and pressure-dependent critical nucleus diameter. If the latter
12 is the case, the inception rate is proportional to the product of the concentra-
13 tions of the collision partners and the transition regime coagulation kernel.
14 More details can be found in [1] and [16].

15
16
17
18
19
20
21 *Condensation:* An existing particle can grow through (barrier-free) de-
22 position of SiH_2 , $\text{Si}_2\text{H}_4\text{A}$, or $\text{Si}_2\text{H}_4\text{B}$ molecules from the gas phase onto its
23 surface. It is assumed that the collision efficiency, *i.e.* the probability of
24 sticking, is unity. The rate is given by a free-molecular collision kernel.

25
26
27
28
29
30
31 *Surface reaction:* Apart from simply condensing, gas-phase species can
32 also react heterogeneously on the particle surface. Specifically, silanes (SiH_4 ,
33 Si_2H_6 , and Si_3H_8) can be integrated into the particle, with each step releasing
34 one, two, and three molecules of hydrogen, respectively. The rate is propor-
35 tional to the particle surface area and an Arrhenius expression with non-zero
36 activation energy. Rounding of adjacent primary particles caused by this
37 process is also taken into account.

38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53 *Hydrogen release:* In order to attain a stable crystal structure, particles
54 need to release some of the hydrogen acquired through each of the above
55 processes. The rate of desorption is proportional to an Arrhenius expression
56 and the coverage of hydrogen on the particle surface, which is approximated
57 by the ratio of hydrogen to silicon atoms within the particle. It is assumed
58 that the sintering level of adjacent primaries is unaffected by this process,
59 *i.e.* the connectivity matrix remains unchanged.

60
61
62
63
64
65
Coagulation: Two particles can collide and stick to each other at their

point of contact. The rate is given by transition regime coagulation kernel, which is the harmonic mean of the slip-flow and free-molecular kernels. The transition kernel is valid across a wide range of Knudsen numbers, and thus wide ranges of pressures and particle sizes (see [19] and [22] for more details).

Sintering: The sintering of any pair of adjacent primary particles is modelled by an exponential decay of the excess of the joint surface area of the primaries compared to the surface area of their equivalent sphere. In other words, the corresponding entry in the connectivity matrix decreases exponentially towards the equivalent spherical area of the primary particle pair.

2.2. Omission-based regression influence diagnostics

2.2.1. Parameter estimation

Given a set of N experimental observations η_n^{exp} , with $n = 1, \dots, N$. For example, these could be, as in this work, means or modes of the particle size distribution at given temperatures and pressures. Assuming we have a model which depends on a vector ϑ of P model parameters, we denote its response for the conditions of the n^{th} experiment by $\eta_n(\vartheta)$. For simplicity, we restrict ourselves in this work to a single response, but the generalisation of all that follows to multiple responses is straightforward.

In order to quantify agreement between experiment and model, a measure of the distance between the model response and experimental results needs to be defined. We use the ordinary least-squares objective function

$$\Phi(\vartheta) := \sum_{n=1}^N [\eta_n(\vartheta) - \eta_n^{\text{exp}}]^2 \quad (1)$$

for this purpose. The term ‘ordinary’ refers to the fact that the covariance matrix of the responses is the unity matrix, *i.e.* the responses are assumed

to be uncorrelated and are subject to the same or very similar uncertainties, meaning all the terms in the sum are equally weighted.

The vector $\hat{\vartheta}$ of parameter values which are optimal with respect to the objective function can be obtained by minimising (1):

$$\hat{\vartheta} := \underset{\vartheta}{\operatorname{argmin}} \Phi(\vartheta) \quad (2)$$

The best estimate of the model responses is then defined as $\hat{\eta} := \eta(\hat{\vartheta})$.

2.2.2. Influence measures

The basic idea underlying omission-based regression influence diagnostics is to analyse the effect of deleting a single observation from the considered set of data. In the following, we use a subscript “ $-i$ ” to denote quantities based on the data set with the i^{th} observation removed. In particular, the objective function (1) becomes

$$\Phi_{-i}(\vartheta) := \sum_{n=1, \dots, i-1, i+1, \dots, N} [\eta_n(\vartheta) - \eta_n^{\text{exp}}]^2, \quad (3)$$

with the corresponding best parameter estimate

$$\hat{\vartheta}_{-i} := \underset{\vartheta}{\operatorname{argmin}} \Phi_{-i}(\vartheta) \quad (4)$$

and response estimate $\hat{\eta}_{-i} := \eta(\hat{\vartheta}_{-i})$.

There are numerous ways of assessing how the optimum, *i.e.* the best estimate of the parameters, is affected by removing a data point [7]. The most elementary statistic is obtained by considering the difference between the best estimate of the parameters and the best estimate with the i^{th} data point removed:

$$D_{ij}^* := \hat{\vartheta}_j - \hat{\vartheta}_{-i,j}, \quad (5)$$

where $\hat{\vartheta}_{-i,j}$ is the value of the j^{th} parameter obtained from the optimisation with the i^{th} experiment omitted. In the literature this is usually referred to as DFBETA_{*i*} [23, p. 13].

We note that such an analysis requires $\hat{\vartheta}_{-i}$ to be calculated for all $i = 1, \dots, N$, each requiring one optimisation. This can become computationally prohibitively expensive if the model itself is expensive or there are many experimental observations. If the considered model is linear, at least approximately, then it is possible to derive a formula which allows calculating the entire set of D_{ij}^* based on only a single optimisation [2]. This, however, is not an option if the model responses are strongly non-linear or are subject to numerical or statistical noise. The model considered in this work is by nature a stochastic model and its responses do exhibit non-negligible noise.

In order to compare or rank different parameters against each other with respect to their influence, due to different physical dimensions and/or orders of magnitude, it is essential to consider non-dimensionalised diagnostic measures. Belsley et al. [23, p. 13] recommend to normalise by the square root of an estimate of the variance of each parameter (with the i^{th} data point removed). This allows assessing the influence of data points on each parameter in relation to their uncertainty. Specifically, they propose to measure the influence of the i^{th} experiment upon the j^{th} parameter using $\text{DFBETAS}_{ij} := D_{ij}^*/(\text{Var } \hat{\vartheta}_j)^{1/2}$ (see also [24]), where $\text{Var } \hat{\vartheta}_j$ refers to the variance of the j^{th} parameter. In some situations, the parameter variance may not be readily available, such as in this work where we directly optimise the model while progressively excluding experiments. Hence, we simply use here parameters which are normalised by (logarithmically) mapping them to the

interval $[-1, 1]$.

Cook's distance [5], one of the most widely-used influence diagnostics, can be a useful tool for assessing the influence of an experimental data point during an optimisation. In the special case we consider in this work, *i.e.* that of uncorrelated responses with similar uncertainty, it can be defined as [7]

$$C_i := \frac{\sum_{n=1}^N [\hat{\eta}_n - \hat{\eta}_{-i,n}]^2}{Ps^2}, \quad (6)$$

where $\hat{\eta}_{-i,n}$ is the value of the model response for the conditions of the n^{th} experiment obtained using the best parameter value estimates determined through optimisation with the i^{th} observation omitted (*i.e.* $\hat{\vartheta}_{-i}$), and where s^2 is an estimate of the mean square error, given by

$$s^2 = \frac{1}{N-P} \sum_{n=1}^N (\eta_n^{\text{exp}} - \hat{\eta}_n)^2. \quad (7)$$

Large values of Cook's distance C_i occur if deleting case i causes large differences in the parameter estimates.

The motivation for definition (6) stems from the notion of joint confidence regions for the parameters. Joint $100(1 - \alpha)\%$ confidence ellipsoids for the model responses can be defined as

$$(\hat{\eta} - \eta)^\top \Sigma^{-1} (\hat{\eta} - \eta) \leq Ps^2 F(P, N - P, 1 - \alpha), \quad (8)$$

with s given by (7), and $F(P, N - P, 1 - \alpha)$ the $1 - \alpha$ point of the F -distribution (consult [25, pp. 94 & 108] and [26] for more details). Σ is the covariance matrix of the responses. Cook introduced his original measure [5, 27] for ordinary least squares, *i.e.* unity covariance matrix, and later generalised it to weighted least squares [3, p. 209]. As mentioned above, if the responses

are uncorrelated, of equal dimension, and of similar order of magnitude and uncertainty, Σ can be assumed to be the unit matrix.

Definition (6), like (5), involves one optimisation per experimental data point. As mentioned above, in situations where this is too computationally expensive, there may be the option of conducting a linearised analysis. For linear models, one can derive an expression for Cook's distance (6) which requires only a single regression for all observations. Whether or not a linear approximation is appropriate can be decided for example by means of local curvature [28, 29], but this is beyond the scope of the paper.

It is noted that Cook's distance measures only the *overall* influence of an observation, in contrast to (5), which assesses parameters individually. More generally, while in this work we consider the influence of single observations only on either single parameters or the model as a whole, this can be generalised to the influence of subsets of observations on subsets of parameters in the model (see for example [24, 7]). As the original notions, however, the measures tend to be applicable to linear models only, and may require additional regressions.

It is furthermore noted that, unlike (5), the Cook distance (6) is dimensionless by definition – a necessary property in order to achieve a generic classification of data points.

In a wider context, recall that a more traditional way to examine the influence of a data point on parameter estimates would be to conduct a sensitivity analysis of the best estimates with respect to the measured data [12], also known as perturbation of the optimum [30] (see also [2]). That is, consider $\partial\hat{\theta}_j/\partial\eta_n^{\text{exp}}$, with (2) and (1). However, approximating such derivatives

by finite differences is problematic for stochastic, or otherwise noisy models, such as in this application, as mentioned above. Additionally, omitting a data point is not local in the sense that it causes a finite step-change in the objective function rather than a small continuous change resulting from a small perturbation of the data point, as is implied by the use of derivatives in a sensitivity analysis.

2.2.3. Outlier detection

One way of identifying potential outliers is by means of a threshold: A data point is deemed to require further attention if the corresponding value of the chosen diagnostic measure exceeds the threshold. Naturally, the choice of any such threshold is ultimately arbitrary, which is reflected in the fact that a range of them has been suggested in the literature. For example, Bollen and Jackman [31] propose

$$C_i \geq 4/N. \quad (9)$$

This threshold is very conservative in the sense that it tends to highlight too many points as outliers. On the other hand, Cook and Weisberg [32, p. 345] suggest

$$C_i \geq 1, \quad (10)$$

i.e. approximately the median of the F distribution with P and $N - P$ degrees of freedom (see Eqn. (8)). Irrespective of which value is chosen, it needs to be emphasised that this method can give only a rough indication, which should be interpreted merely as a suggestion of which data points warrant closer investigation. The main reason for this is that the method does not automatically distinguish between errors and highly influential points

which potentially point towards genuine model improvements. Therefore, highlighted points should not necessarily be excluded from the analysis, as one may lose valuable information. Furthermore, whether or not a data point is deemed an ‘outlier’ by this method, is by definition dependent on the chosen model. That is, a data point labelled an outlier with respect to one model, may or may not appear as an outlier with respect to another (possibly better) model. As there is no consensus in the literature as to which cut-off should be used, in this work we consider both (9) and (10).

3. Experimental data

As in previous work [1, 20], a total of nineteen experimental data points were selected from six different studies, spanning a range of process conditions and reactor configurations. Reactor types include hot-wall flow reactors and a shock tube, for each of which different temperatures, pressures, residence times, and initial silane mole fractions are covered. The particular selection of studies, an overview of which is given in Table 2, was motivated by covering a range of conditions. This choice is, however, arbitrary amongst large amounts of literature (too much to review here comprehensively), which include further works on hot-wall reactors [39, 40, 41, 42], microwave reactors [43, 44, 45], and plasma reactors [46, 47], to name but a few.

The study of Körner et al. [33] is focused on synthesising silicon nanoparticles with narrow size distributions in a hot-wall flow reactor. In this setup, it turns out that most of the precursor is lost to deposits on the reactor wall, and therefore the initial composition is adjusted to account for this particle deposition [48]. As in [49], an initial silane mass of about 6×10^{-5} kg/m³ is

Table 2: Experimental data sets with process conditions used to model them. X_{SiH_4} denotes the initial silane mole fraction, and τ denotes the residence time.

Idx. i	Reference	Reactor type	Bath gas	X_{SiH_4} [%]	T [K]	P [kPa]	τ [ms]	d type	μ type	μ_i^{exp} [nm]
1	Körmer <i>et al.</i> [33]	Hot-wall flow reactor	Ar	4.0	873-1373	2.5	80	d_{pri}	Mode	26.7
2				4.0	873-1373		192		Mean	26.0
3				12.0	873-1373		192		Mean	38.0
4				12.8	873-1373		80		Mode	31.0
5				2.0	873-1373		80		Mode	41.0
6				8.0	873-1373		80		Mode	24.0
7				4.0	873-1373		420		Mode	32.5
8				4.0	673-1173		420		Mode	21.2
9				4.0	773-1273		420		Mode	28.5
10	Frenklach <i>et al.</i> [34]	Shock tube	Ar	3.3	1089	49	2.6	d_{pri}	Mode	11.0
11					1320		2.1			11.0
12					1580		1.8			15.0
13	Wu <i>et al.</i> [35]	Hot-wall flow reactor	N ₂	1.0	770-1520	101	1000	d_{mob}	Mode	127
14	Flint <i>et al.</i> [36]	Laser-	Ar	21.4	923-1270	20	5.2	d_{pri}	Mean	43.4
15		driven		9.0	1023-1483		18			55.4
16		flow reactor		0.6	1023-1400		53			23.0
17	Nguyen and Flagan [37]	Hot-wall	N ₂	0.1	770-1080	101	900	d_{mob}	Mode	89.0
18		flow reactor		0.04						51.0
19	Onischuk <i>et al.</i> [38]	Hot-wall flow reactor	Ar	5.0	853	39	870	d_{pri}	Mean	52.0

assumed. The amount of mass expected for a partial pressure of 1 mbar of silane at 1024 K is about 3.8×10^{-4} kg/m³ indicating that only about 16% of the precursor are available to form particles. The initial silane fractions listed in Table 2 for this data subset are adjusted accordingly for our simulations.

The Flint *et al.* [36] data refers to their cases 630S, 631S, and 654S, respectively. The experiment is described in detail in [50, 51, 52], including how to convert flow rates into residence times and initial compositions.

Further work from the same group includes [53, 54].

4. Results and discussion

Both reactor types occurring in the set of experiments (Table 2), *i.e.* flow reactor and shock tube, are modelled as homogenous batch reactors. The shock tube is modelled as a constant temperature, constant pressure reactor. For the flow reactors, plug-flow is assumed, and the experimentally measured temperature profile, where available, is imposed. In case 19 [38], no temperature profile is available, so a constant temperature is assumed, and the residence time given refers to the approximate time spent in the ‘hot-zone’, *i.e.* at that temperature.

As software to carry out the necessary optimisations, we use the Model Development Suite (MoDS) [55] – a software tool for conducting various generic tasks to develop black-box models. Such tasks include parameter estimation and uncertainty quantification [56], Design of Experiments (DoE) [57], and global sensitivity analysis [20].

Each optimisation involved in the Cook distance and DFBETA analysis is performed in two stages: Firstly, a quasi-random global search is conducted using Sobol low-discrepancy sequences [58]. Secondly, starting from the best point identified in the first stage, a local optimisation is carried out using the Simultaneous Perturbation Stochastic Approximation (SPSA) [59, 60] algorithm. The SPSA method estimates the local gradient based on only two objective function evaluations, and can be shown to obey the traditional gradient descent *on average*. It is designed for problems where stochastic noise is present. The motivation for the first stage is to avoid becoming

trapped in local minima or valleys on the objective function surface, which could happen with a method purely based on the local gradient. Chemical kinetic objective functions are widely reported to exhibit a complex, highly structured surface with multiple local minima and/or valleys (see for example [26]). Regarding the second stage, the reason for not choosing a more conventional method utilising the local Jacobi matrix or Hessian is the stochastic noise in the model response. While the procedure adopted here cannot guarantee to find the global minimum, based on previous experience [56], a low-lying minimum can be found at a manageable computational expense. On an objective function surface with multiple local minima, there is then of course the risk of selecting the ‘wrong’ optimum, *i.e.* not the global one. Any conclusions derived from perturbations such as those induced by omission of data points may change depending on the chosen minimum and the local geometry surrounding it.

Table 3: The seven model parameters considered in the influence analysis, all Arrhenius pre-exponential factors (see Table 1), with optimal values resulting from optimisation against the complete data set.

Idx.	Param.	Opt. value	Unit	Phase	Role
1	$A_{1,LP}$	2.87×10^{12}			Low-pressure limit of reaction #1
2	$A_{2,LP}$	2.11×10^{35}			Low-pressure limit of reaction #2
3	$A_{3,LP}$	4.90×10^{39}	$\text{cm}^3/\text{mol/s}$	Gas	Low-pressure limit of reaction #3
4	$A_{5,LP}$	2.98×10^{68}			Low-pressure limit of reaction #5
5	$A_{8,rev}$	1.48×10^{14}			Reverse of reaction #8
6	A_{SR, SiH_4}	4.47×10^{33}	$\text{cm}/\text{mol/s}$	Particle	Surface reac.: silane addition, H ₂ -release
7	A_{H_2}	1.88×10^{18}	$1/\text{s}$		H ₂ -release from particle

Here, seven parameters were adjusted which represent key gas-phase and heterogeneous growth rates identified through sensitivity analysis [1]. We

note that this choice is consistent with reports in the literature [15] which suggest that it is the low-pressure limit that is of interest to the conditions considered in this work. Details are given in Table 3. Thus, the vector of model parameters to be optimised is given by

$$\vartheta = (A_{1,\text{LP}}, A_{2,\text{LP}}, A_{3,\text{LP}}, A_{5,\text{LP}}, A_{8,\text{rev}}, A_{\text{SR},\text{SiH}_4}, A_{\text{H}_2}).$$

The optimal values for the parameters resulting from optimisation against the full data set are also given in Table 3. The differences in these values as compared to [1] and [20] are due to the fact that different sets of responses are being considered.

For the optimisation against the complete data set, 800 Sobol points were generated, followed by 240 SPSA points. Recall that each point involves one evaluation of the objective function (Eqn. 1), and that every objective function evaluation involves 19 model evaluations. For all subsequent optimisations, *i.e.* those of Φ_{-i} (Eqn. 3) with $i = 1, \dots, 19$, the model evaluations performed as part of the original set of Sobol points can be re-used, as all that is required is for each i to calculate the different objective function Φ_{-i} for all of the points. For each of the Φ_{-i} optimisations, 120 SPSA points were used. In total, this corresponds to about 3300 CPU-hours of computation.

The Cook distance analysis was conducted for all of the 19 experiments in Table 2, and results are shown in Fig. 1. In this figure, the responses are grouped by the particular experimental papers from which they were obtained. Both of the two outlier thresholds, Eqn. (9) and Eqn. (10), are shown. While several of the observations exceed the lower threshold (9), only two of them exceed the upper one (10) (with one of them only marginally). This is consistent with reports that (9) is too conservative in that it has a ten-

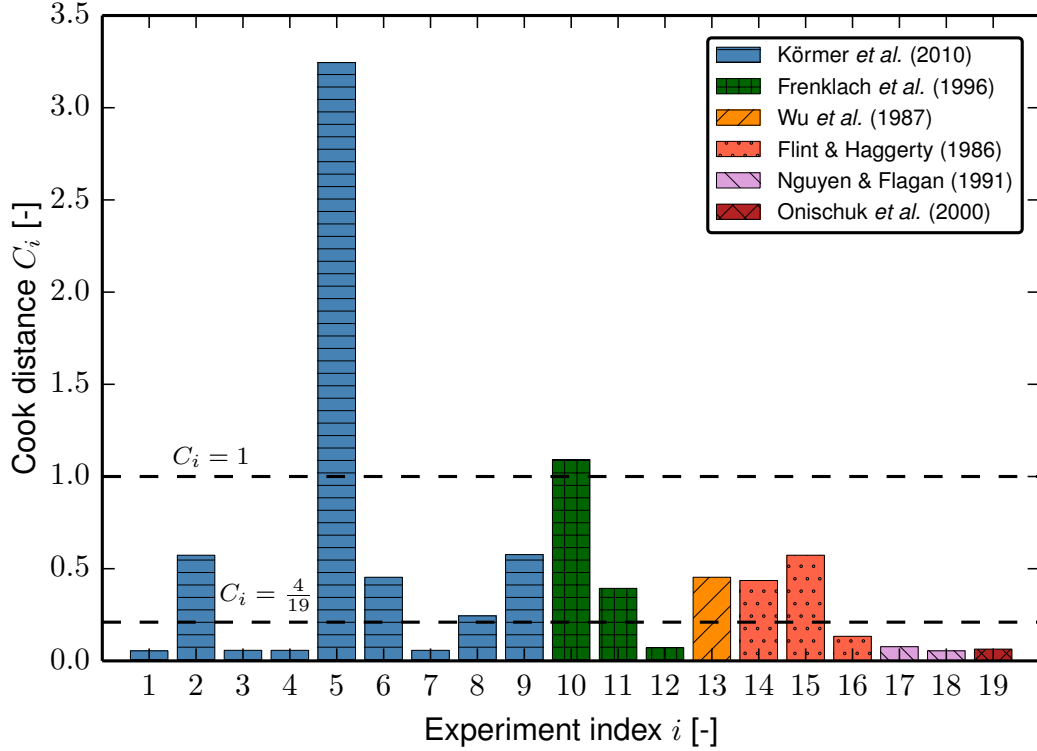
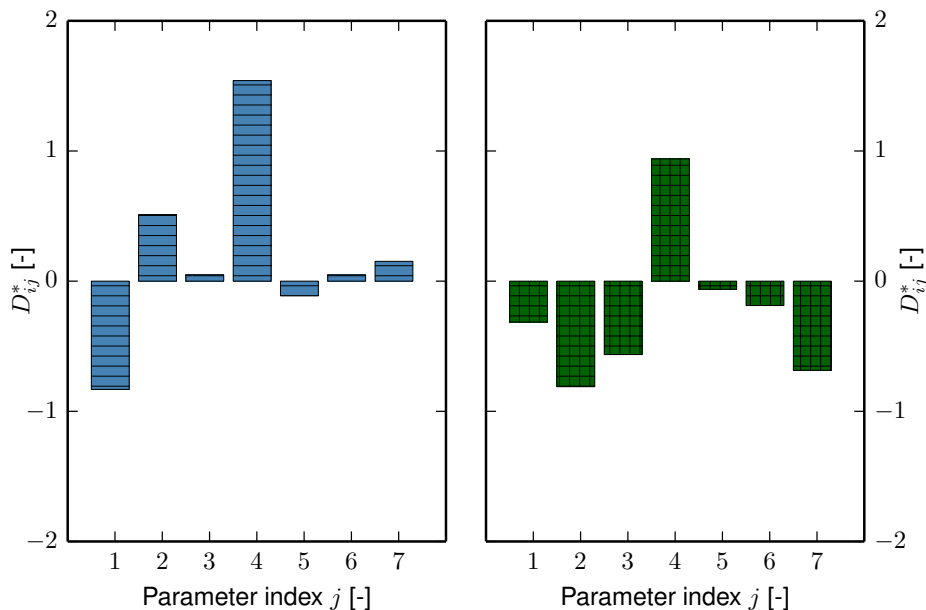


Figure 1: Overall influence of each of the experimental observations in Table 2 as measured by Cook's distance C_i (Eqn. 6). Each of the thresholds (9) and (10) are indicated through dashed horizontal lines.

dency to highlight too many observations, as mentioned in subsection 2.2.3. We conclude that observation $i = 5$ requires further attention, as its Cook distance exceeds both thresholds and is significantly larger than all the others. This indicates that this experimental point most strongly affects the objective function Φ (Eqn. 1), which in turn affects the parameter estimates, *i.e.* the optimal values $\hat{\vartheta}$ of the parameters (Eqn. 2). It could furthermore suggest that this particular observation might be an outlier with respect to the present model, or, more likely, that the model describes it inadequately.



(a) Influence of observation $i = 5$ by Körner et al. [33] on each of the considered model parameters. (b) Influence of observation $i = 10$ by Frenklach et al. [34] on each of the considered model parameters.

Figure 2: DFBETA D_{ij}^* (Eqn. 5), for the two most influential experimental observations as identified in Fig. 1 (see also Table 2), for each of the parameters in Table 3.

Additionally, a DFBETA analysis was conducted to assess how the experimental observations affect the values of the parameters which are determined through the optimisation (Fig. 2). In terms of highlighting individual observations, the DFBETA analysis agrees with the Cook distance analysis: The values of D_{ij}^* for $i = 5$ and $i = 10$ are at least two orders of magnitude larger than those obtained for any other experiment. The DFBETA values for these experiments are shown in Figs. 2a and 2b respectively. We notice that the best estimate of parameter 4, *i.e.* the pre-exponential factor in the low-pressure limit of reaction 5 (Table 1), is influenced most by both of the

considered experimental observations.

In principle, there are two possible reasons for why an observation stands out in a Cook distance or DFBETA analysis: errors associated with the experiments, and errors associated with the model. Regarding experimental errors, we assume here that all experimental data are both correct and accurate. Considering model errors, these can be further categorised into the following: errors arising from the solution methodology, *i.e.* numerical algorithms, and flaws in the model. Specifically in this case, the latter include reactor model errors, and deficiencies in the gas or particulate phase sub-models.

Figure 3 shows particle size distributions for those experiments in Table 2 for which they have been measured. Two sets of model results are shown – one for the optimisation against the complete data set, and one for the data set with the 5th experiment omitted. As expected, if the 5th experiment is omitted, the corresponding response deteriorates significantly (Fig. 3e). Recall that only the means or modes of the distributions are optimised, not the widths or any other characteristic. This is most obvious in cases 7-9 (Figs. 3g-i) for example, where the modes agree reasonably well but the model distributions are noticeably wider than the experimental ones. Note, however, that adding, say, the standard deviation of the distributions as optimisation targets by itself, *i.e.* without adding further degrees of freedom in terms of model parameters to be optimised, will not ‘improve’ the fit in any way. The fit can only improve if model parameters are included in the optimisation which are suitable in the sense that they affect the width of the distributions independently of the mean, provided such degrees of freedom

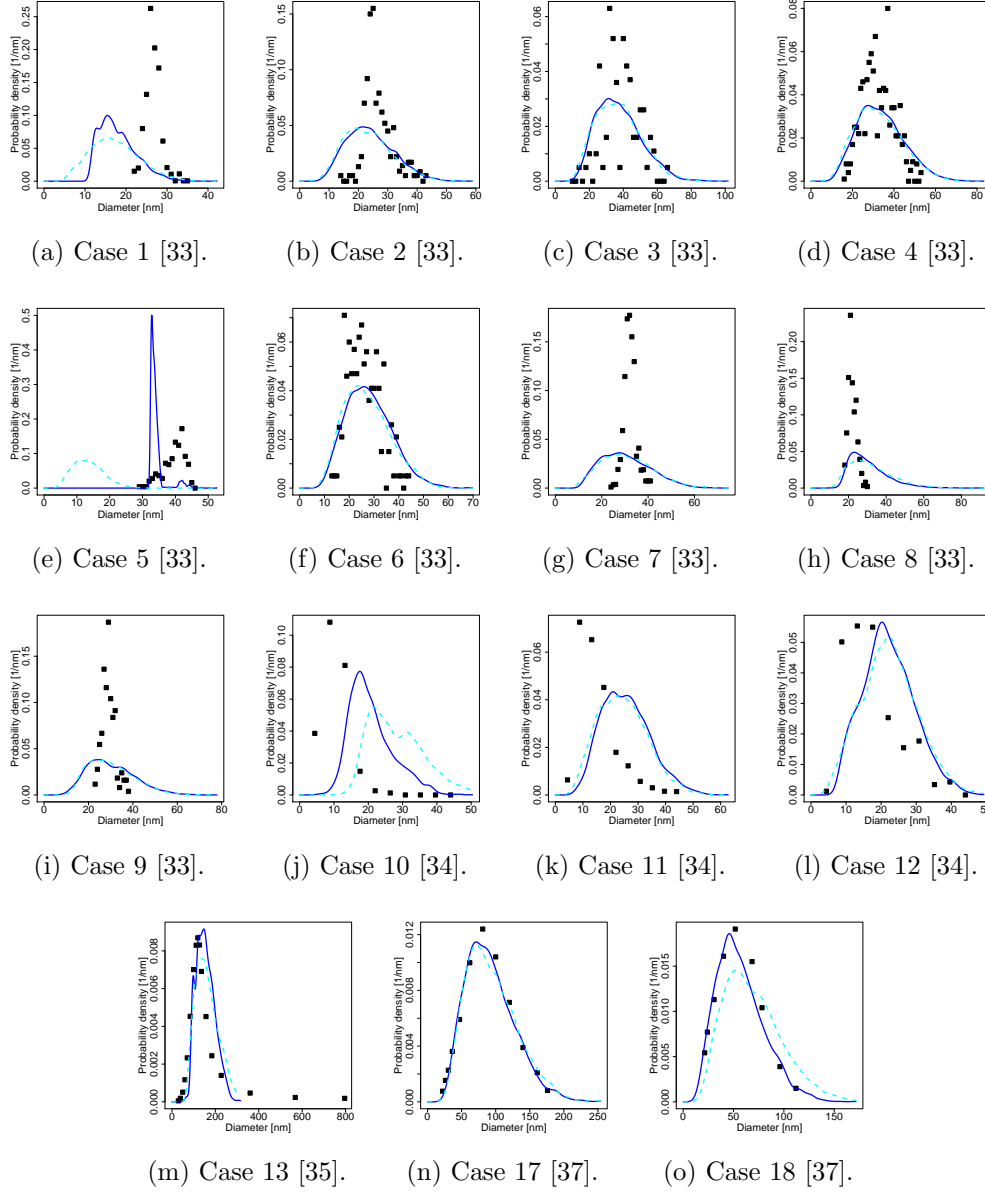


Figure 3: Particle size distributions for those experiments for which they were measured. Solid lines: model optimised against the complete data set. Dashed lines: model optimised against the data set with the 5th experiment omitted. Points: experiment.

exist in the model.

The $i = 5$ experiment refers to the case of lowest initial silane concentration (0.5 mbar partial pressure) reported by Körner et al. [33]. In this hot-wall reactor experiment, a modal size of 41 nm was obtained for the primary particles, larger than that obtained for higher concentrations: 1 mbar partial pressure yielded 27 nm primaries, 2 mbar yielded 24 nm primaries, rising again to 31 nm at 4 mbar, all at the same residence time (and total pressure). This inverse relationship for the smaller initial concentrations is not captured by the model, thus indicating that this aspect requires further development. More specifically, this suggests that the ratio between the inception and condensation rates should be revisited, as this directly controls the size and number of primary particles. Given that the inception mechanism in particular remains an active area of research, with several fundamental open questions, this might be the most natural starting point.

In order to investigate further the kinetic role of initial silane concentration and total pressure, we conducted flux analyses of Si, time-integrated as well as instantaneous, for a range of concentrations and pressures, covering the conditions of all experiments in Table 2. At high dilutions, the importance of unimolecular reactions is expected to increase relative to bimolecular ones. However, we found that, for the mechanism used, whilst the pressure dependence of the net fluxes can be significant, their dependence on dilution is relatively minor within the range considered. Besides, we note that even though experiment 5 is the most dilute amongst the ones by Körner et al. [33], experiments 13, 16, 17, and 18 are more dilute, some significantly so (case 18 by a factor of 50), with the model agreeing well especially with

the latter three. Therefore, irrespective of the performance of the gas-phase mechanism at low silane concentrations, this alone is insufficient to explain the overall model behaviour for experiment 5.

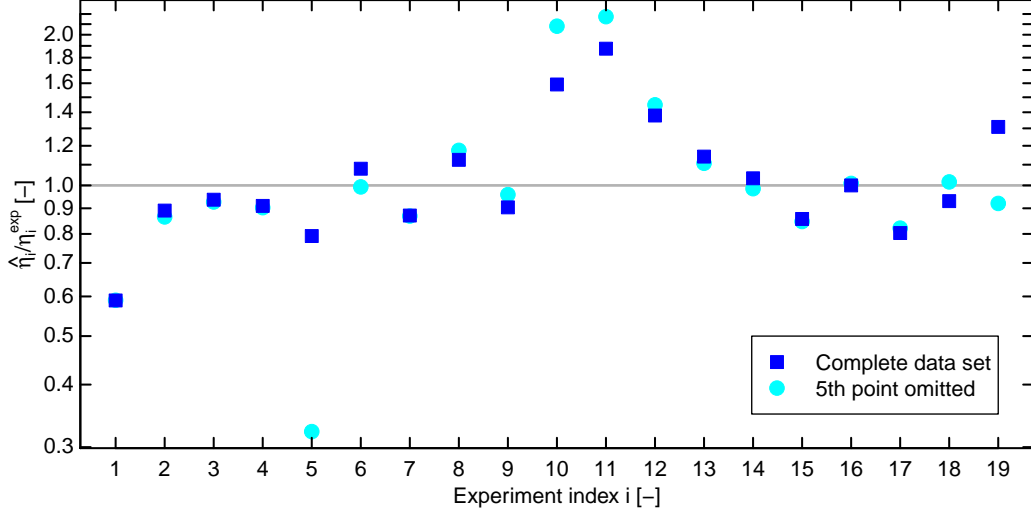


Figure 4: Ratios of model responses to experimental values for each of the 19 experiments in Table 2. Squares: model optimised against the complete data set. Circles: model optimised against the data set with the 5th experiment omitted.

Figure 4 shows ratios of model responses to experimental ones for all experiments. Again, two sets of results are shown – one for the optimisation against the complete data set, and one for the data set with the 5th experiment omitted, and the most obvious feature is again the worsening of the response corresponding to the 5th experiment. Even though some responses, such as for cases 10 and 11 (Figs. 3j and k), have deteriorated, one should note that the value of the overall objective function (Eqn. 3) is still lower for the omitted set than the full set. This is mainly due to responses 18 and 19, and also 13, improving and their absolute values being much larger

than those of responses 10 and 11. Referring again to Table 2, this hints at a competition or trade-off between two very different scenarios which the model is not capable of capturing simultaneously: the short residence-time regime with early, nucleation-stage particles, versus the longer residence-time regime with mature, larger aggregates.

5. Conclusions

We determined optimal values of seven parameters in a population balance model for the formation of silicon nanoparticles by means of least-squares optimisation against a set of 19 experiments. The influence of each of those measurements on the values of the considered kinetic model parameters was then quantified using Cook’s distance and DFBETA – two basic omission-based measures popular in the field of regression influence diagnostics. An outlier analysis was then conducted by applying standard thresholds in order to identify the most important experimental datasets in the optimisation. This highlighted one particular experimental condition for further scrutiny. We emphasise again that, in general, a particular measurement exceeding an outlier threshold does not necessarily imply that there is a problem with that measurement or more generally the experiment. In the first instance, one should thoroughly examine whether there are shortcomings in the model which are responsible for the disagreement with the measurement. This informs future model development [57] by helping to identify aspects of the model which require improvement. Furthermore, if one regards the model as a formal representation of the best current knowledge about the experiment or system under consideration [61], then the methods can be thought of as

giving an indication as to which measurements are most informative.

Acknowledgements

This work was partly funded by the Cambridge Australia Trust, by the National Research Foundation (NRF), Prime Minister’s Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CRE-ATE) programme, and by the European Union Horizon 2020 Research and Innovation Programme under grant agreement 646121.

References

- [1] W. J. Menz, M. Kraft, A new model for silicon nanoparticle synthesis, *Combust. Flame* 160 (2013) 947–958. doi:10.1016/j.combustflame.2013.01.014.
- [2] S. Mosbach, M. Kraft, Influence of experimental observations on n-propylbenzene kinetic parameter estimates, *Proc. Combust. Inst.* 35 (2015) 357–365. doi:10.1016/j.proci.2014.05.061.
- [3] R. D. Cook, S. Weisberg, *Residuals and Influence in Regression*, Chapman and Hall, New York, 1982.
- [4] R. Schall, T. T. Dunne, Influential variables in linear regression, *Technometrics* 32 (1990) 323–330. doi:10.1080/00401706.1990.10484685.
- [5] R. D. Cook, Detection of influential observation in linear regression, *Technometrics* 19 (1977) 15–18.

- [6] N. R. Draper, J. A. John, Influential observations and outliers in regression, *Technometrics* 23 (1981) 21–26. doi:10.1080/00401706.1981.10486232.
- [7] S. Chatterjee, A. S. Hadi, Influential observations, high leverage points, and outliers in linear regression, *Stat. Sci.* 1 (1986) 379–393. doi:10.1214/ss/1177013622.
- [8] A. S. Tomlin, The role of sensitivity and uncertainty analysis in combustion modelling, *Proc. Combust. Inst.* 34 (2013) 159–176. doi:10.1016/j.proci.2012.07.043.
- [9] R. Feeley, P. Seiler, A. Packard, M. Frenklach, Consistency of a reaction dataset, *J. Phys. Chem. A* 108 (2004) 9573–9583. doi:10.1021/jp047524w.
- [10] A. V. Fiacco, G. P. McCormick, *Nonlinear Programming – Sequential Unconstrained Minimization Techniques*, Classics in Applied Mathematics, SIAM, 1990.
- [11] L. Eno, J. G. B. Beumee, H. Rabitz, Sensitivity analysis of experimental data, *Appl. Math. Comput.* 16 (1985) 153–163. doi:10.1016/0096-3003(85)90005-0.
- [12] H. Rabitz, M. Kramer, D. Dacol, Sensitivity analysis in chemical kinetics, *Annu. Rev. Phys. Chem.* 34 (1983) 419–461. doi:10.1146/annurev.pc.34.100183.002223.
- [13] T. Turányi, Sensitivity analysis of complex kinetic systems.

Tools and applications, *J. Math. Chem.* 5 (1990) 203–248.
doi:10.1007/BF01166355.

- [14] P. Ho, M. E. Coltrin, W. G. Breiland, Laser-induced fluorescence measurements and kinetic analysis of Si atom formation in a rotating disk chemical vapor deposition reactor, *J. Phys. Chem.* 98 (1994) 10138–10147. doi:10.1021/j100091a032.
- [15] E. L. Petersen, M. W. Crofton, Measurements of high-temperature silane pyrolysis using SiH_4 IR emission and SiH_2 laser absorption, *J. Phys. Chem. A* 107 (2003) 10988–10995. doi:10.1021/jp0302663.
- [16] W. J. Menz, S. Shekar, G. P. E. Brownbridge, S. Mosbach, R. Körner, W. Peukert, M. Kraft, Synthesis of silicon nanoparticles with a narrow size distribution: A theoretical study, *J. Aerosol Sci.* 44 (2012) 46–61. doi:10.1016/j.jaerosci.2011.10.005.
- [17] S. Shekar, A. J. Smith, W. J. Menz, M. Sander, M. Kraft, A multidimensional population balance model to describe the aerosol synthesis of silica nanoparticles, *J. Aerosol Sci.* 44 (2012) 83–98. doi:10.1016/j.jaerosci.2011.09.004.
- [18] S. Shekar, W. J. Menz, A. J. Smith, M. Kraft, W. Wagner, On a multivariate population balance model to describe the structure and composition of silica nanoparticles, *Comput. Chem. Eng.* 43 (2012) 130–147. doi:10.1016/j.compchemeng.2012.04.010.
- [19] W. J. Menz, R. I. A. Patterson, W. Wagner, M. Kraft, Application of stochastic weighted algorithms to a multidimensional

silica particle model, *J. Comput. Phys.* 248 (2013) 221–234.
doi:10.1016/j.jcp.2013.04.010.

- [20] W. J. Menz, G. P. E. Brownbridge, M. Kraft, Global sensitivity analysis of a model for silicon nanoparticle synthesis, *J. Aerosol Sci.* 76 (2014) 188–199. doi:10.1016/j.jaerosci.2014.06.011.
- [21] E. K. Y. Yapp, D. Chen, J. W. J. Akroyd, S. Mosbach, M. Kraft, J. Camacho, H. Wang, Numerical simulation and parametric sensitivity study of particle size distributions in a burner-stabilised stagnation flame, *Combust. Flame* 162 (2015) 2569–2581. doi:10.1016/j.combustflame.2015.03.006.
- [22] A. Kazakov, M. Frenklach, Dynamic modeling of soot particle coagulation and aggregation: Implementation with the method of moments and application to high-pressure laminar premixed flames, *Combust. Flame* 114 (1998) 484–501. doi:10.1016/S0010-2180(97)00322-2.
- [23] D. A. Belsley, E. Kuh, R. E. Welsch, *Regression Diagnostics : Identifying Influential Data and Sources of Collinearity*, John Wiley & Sons, 1980.
- [24] R. D. Cook, S. Weisberg, Characterizations of an empirical influence function for detecting influential cases in regression, *Technometrics* 22 (1980) 495–508. doi:10.1080/00401706.1980.10486199.
- [25] N. R. Draper, H. Smith, *Applied Regression Analysis*, 2nd ed., John Wiley & Sons, New York, 1981.
- [26] M. Frenklach, in: W. C. Gardiner (Ed.), *Combustion Chemistry*, Springer Verlag, New York, 1984, pp. 423–453.

- [27] R. D. Cook, Influential observations in linear regression, *J. Am. Stat. Assoc.* 74 (1979) 169–174.
- [28] D. M. Bates, D. G. Watts, Relative curvature measures of nonlinearity, *J. Roy. Stat. Soc. B* 42 (1980) 1–25.
- [29] G. A. F. Seber, C. J. Wild, *Nonlinear Regression*, John Wiley & Sons, 2003.
- [30] A. V. Fiacco, *Introduction to Sensitivity and Stability Analysis in Nonlinear Programming*, volume 165 of *Mathematics in Science and Engineering*, Academic Press, New York, 1983.
- [31] K. A. Bollen, R. W. Jackman, in: J. Fox, J. S. Long (Eds.), *Modern Methods of Data Analysis*, Sage Publications, Newbury Park, 1990, pp. 257–291.
- [32] R. D. Cook, S. Weisberg, in: S. Leinhardt (Ed.), *Sociological Methodology*, Jossey-Bass, San Francisco, 1982, pp. 313–316.
- [33] R. Körner, M. P. M. Jank, H. Ryssel, H.-J. Schmid, W. Peukert, Aerosol synthesis of silicon nanoparticles with narrow size distribution – Part 1: Experimental investigations, *J. Aerosol Sci.* 41 (2010) 998–1007. doi:10.1016/j.jaerosci.2010.05.007.
- [34] M. Frenklach, L. Ting, H. Wang, M. J. Rabinowitz, Silicon particle formation in pyrolysis of silane and disilane, *Israel J. Chem.* 36 (1996) 293–303. doi:10.1002/ijch.199600041.

- [35] J. J. Wu, H. V. Nguyen, R. C. Flagan, A method for the synthesis of sub-micron particles, *Langmuir* 3 (1987) 266–271. doi:10.1021/la00074a021.
- [36] J. H. Flint, R. A. Marra, J. S. Haggerty, Powder temperature, size, and number density in laser-driven reactions, *Aerosol Sci. Tech.* 5 (1986) 249–260. doi:10.1080/02786828608959091.
- [37] H. V. Nguyen, R. C. Flagan, Particle formation and growth in single-stage aerosol reactors, *Langmuir* 7 (1991) 1807–1814. doi:10.1021/la00056a038.
- [38] A. A. Onischuk, A. I. Levykin, V. P. Strunin, K. K. Sabelfeld, V. N. Panfilov, Aggregate formation under homogeneous silane thermal decomposition, *J. Aerosol Sci.* 31 (2000) 1263–1281. doi:10.1016/S0021-8502(00)00031-8.
- [39] H. Wiggers, R. Starke, P. Roth, Silicon particle formation by pyrolysis of silane in a hot wall gasphase reactor, *Chem. Eng. Technol.* 24 (2001) 261–264. doi:10.1002/1521-4125(200103)24:3<261::AID-CEAT261>3.0.CO;2-K.
- [40] A. A. Onischuk, V. P. Strunin, M. A. Ushakova, V. N. Panfilov, Studying of silane thermal decomposition mechanism, *Int. J. Chem. Kinet.* 30 (1998) 99–110. doi:10.1002/(SICI)1097-4601(1998)30:2<99::AID-KIN1>3.0.CO;2-O.
- [41] J. O. Odden, P. K. Egeberg, A. Kjekshus, From monosilane to crystalline silicon, part I: Decomposition of monosilane at 690–830 K and initial

pressures 0.1-6.6 MPa in a free-space reactor, *Solar Energy Materials and Solar Cells* 86 (2005) 165–176. doi:10.1016/j.solmat.2004.07.002.

- [42] J. O. Odden, P. K. Egeberg, A. Kjekshus, From monosilane to crystalline silicon, part II: Kinetic considerations on thermal decomposition of pressurized monosilane, *Int. J. Chem. Kinet.* 38 (2006) 309–321. doi:10.1002/kin.20164.
- [43] B. Giesen, H. Wiggers, A. Kowalik, P. Roth, Formation of Si-nanoparticles in a microwave reactor: Comparison between experiments and modelling, *Nanopart. Res.* 7 (2005) 29–41. doi:10.1007/s11051-005-0316-z.
- [44] J. Knipping, H. Wiggers, B. Rellinghaus, P. Roth, D. Konjhodzic, C. Meier, Synthesis of high purity silicon nanoparticles in a low pressure microwave reactor, *J. Nanosci. Nanotechnol.* 4 (2004) 1039–1044. doi:10.1166/jnn.2004.149.
- [45] A. Gupta, H. Wiggers, Surface chemistry and photoluminescence property of functionalized silicon nanoparticles, *Physica E* 41 (2009) 1010–1014. doi:10.1016/j.physe.2008.08.033.
- [46] Z. Shen, T. Kim, U. Kortshagen, P. H. McMurry, S. A. Campbell, Formation of highly uniform silicon nanoparticles in high density silane plasmas, *J. Appl. Phys.* 94 (2003) 2277–2283. doi:10.1063/1.1591412.
- [47] N. J. Kramer, R. J. Anthony, M. Mamunuru, E. S. Aydil, U. R. Kortshagen, Plasma-induced crystallization of silicon nanoparticles, *J. Phys. D* 47 (2014) 075202. doi:10.1088/0022-3727/47/7/075202.

- [48] R. Körmer, H.-J. Schmid, W. Peukert, Aerosol synthesis of silicon nanoparticles with narrow size distribution – Part 2: Theoretical analysis of the formation mechanism, *J. Aerosol Sci.* 41 (2010) 1008–1019. doi:10.1016/j.jaerosci.2010.08.002.
- [49] M. Gröschel, R. Körmer, M. Walther, G. Leugering, W. Peukert, Process control strategies for the gas phase synthesis of silicon nanoparticles, *Chem. Eng. Sci.* 73 (2012) 181–194. doi:10.1016/j.ces.2012.01.035.
- [50] W. R. Cannon, S. C. Danforth, J. S. Flint, J. S. Haggerty, R. A. Marra, Sinterable ceramic powders from laser-driven reactions: I, Process description and modeling, *J. Am. Ceram. Soc.* 65 (1982) 324–330. doi:10.1111/j.1151-2916.1982.tb10464.x.
- [51] W. R. Cannon, S. C. Danforth, J. S. Haggerty, R. A. Marra, Sinterable ceramic powders from laser-driven reactions: II, Powder characteristics and process variables, *J. Am. Ceram. Soc.* 65 (1982) 330–335. doi:10.1111/j.1151-2916.1982.tb10465.x.
- [52] J. Flint, J. Haggerty, A model for the growth of silicon particles from laser-heated gases, *Aerosol Sci. Tech.* 13 (1990) 72–84. doi:10.1080/02786829008959425.
- [53] M. Meunier, J. H. Flint, J. S. Haggerty, D. Adler, Laser-induced chemical vapor deposition of hydrogenated amorphous silicon. I. Gas-phase process model, *J. Appl. Phys.* 62 (1987) 2812–2821. doi:10.1063/1.339412.

- [54] M. Meunier, J. H. Flint, J. S. Haggerty, D. Adler, Laser-induced chemical vapor deposition of hydrogenated amorphous silicon. II. Film properties, *J. Appl. Phys.* 62 (1987) 2821–2829. doi:10.1063/1.339413.
- [55] cmcl innovations, MoDS (Model Development Suite), version 0.2.3, 2015. <http://www.cmclinnovations.com/mod-suite/>.
- [56] S. Mosbach, J. H. Hong, G. P. E. Brownbridge, M. Kraft, S. Gudiyaella, K. Brezinsky, Bayesian error propagation for a kinetic model of n-propylbenzene oxidation in a shock tube, *Int. J. Chem. Kinet.* 46 (2014) 389–404. doi:10.1002/kin.20855.
- [57] S. Mosbach, A. Braumann, P. L. W. Man, C. A. Kastner, G. P. E. Brownbridge, M. Kraft, Iterative improvement of Bayesian parameter estimates for an engine model by means of experimental design, *Combust. Flame* 159 (2012) 1303–1313. doi:10.1016/j.combustflame.2011.10.019.
- [58] I. M. Sobol, On the systematic search in a hypercube, *SIAM J. Numer. Anal.* 16 (1979) 790–793.
- [59] J. C. Spall, Implementation of the Simultaneous Perturbation Algorithm for stochastic optimization, *IEEE T. Aero. Elec. Sys.* 34 (1998) 817–823. doi:10.1109/7.705889.
- [60] T. Hirokami, Y. Maeda, H. Tsukada, Parameter estimation using simultaneous perturbation stochastic approximation, *Electr. Eng. Jpn* 154 (2006) 30–39. doi:10.1002/eej.20239.

- [61] M. Frenklach, Transforming data into knowledge – Process Informatics for combustion chemistry, *Proc. Combust. Inst.* 31 (2007) 125–140. doi:10.1016/j.proci.2006.08.121.

List of Figures

- 1 Overall influence of each of the experimental observations in Table 2 as measured by Cook's distance C_i (Eqn. 6). Each of the thresholds (9) and (10) are indicated through dashed horizontal lines. 19
- 2 DFBETA D_{ij}^* (Eqn. 5), for the two most influential experimental observations as identified in Fig. 1 (see also Table 2), for each of the parameters in Table 3. 20
- 3 Particle size distributions for those experiments for which they were measured. Solid lines: model optimised against the complete data set. Dashed lines: model optimised against the data set with the 5th experiment omitted. Points: experiment. . . . 22
- 4 Ratios of model responses to experimental values for each of the 19 experiments in Table 2. Squares: model optimised against the complete data set. Circles: model optimised against the data set with the 5th experiment omitted. 24